

# Developing a Knowledge Base for NASA Earth Science and Hydrologic Applications

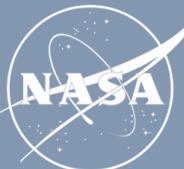
Amanda Weigel<sup>1,2</sup>, Patrick Gatlin<sup>1,3</sup>, Rahul Ramachandran<sup>1,3</sup>, JJ Miller<sup>1,2</sup>, Manil Maskey<sup>1,3</sup>, Jia Zhang<sup>1,3</sup>, Emily Berndt<sup>3</sup>

**NASA/MSFC Data Science Informatics Group<sup>1</sup>**

University of Alabama in Huntsville<sup>2</sup>

NASA Marshall Space Flight Center<sup>3</sup>

Carnegie Mellon University<sup>4</sup>

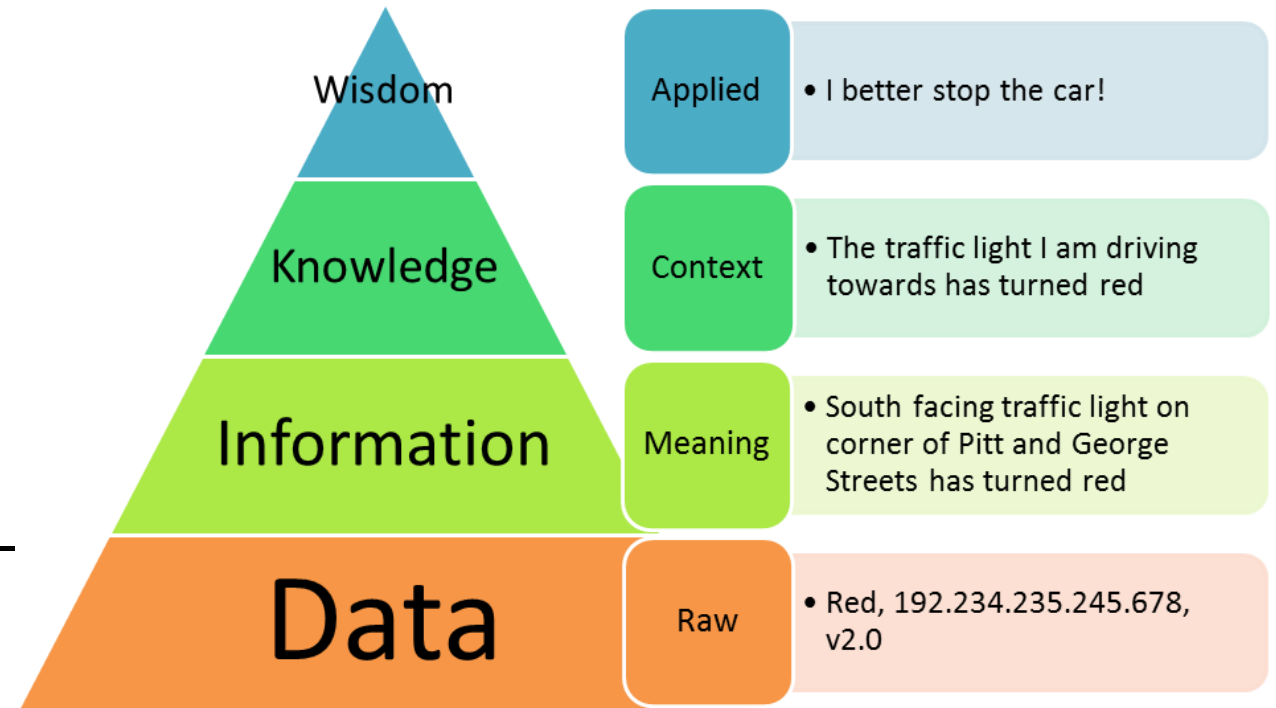




# Data Use Challenges

What common challenges do Earth science and hydrologic data users encounter?

- Data **discovery**
- Data **use**
- Identifying **key resources** about the data.
  - Accessing introductory material (for unfamiliar users).
- Determining what **methods** to use – data processing, quality control and analysis.



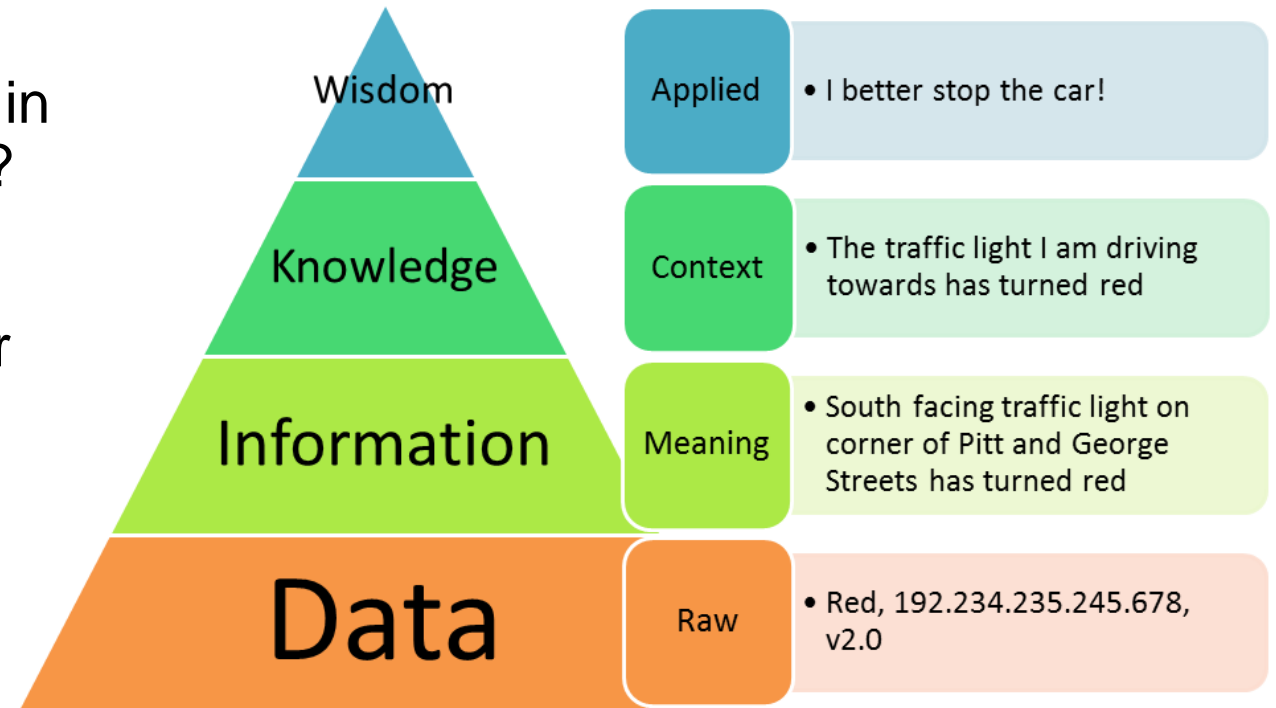
© 2011 Angus McDonald

# Data Use Challenges

**To address these challenges, what difficulties are presented?**

How can data and resources be linked in order to improve the data spin-up time?

How can we work to educate unfamiliar users?



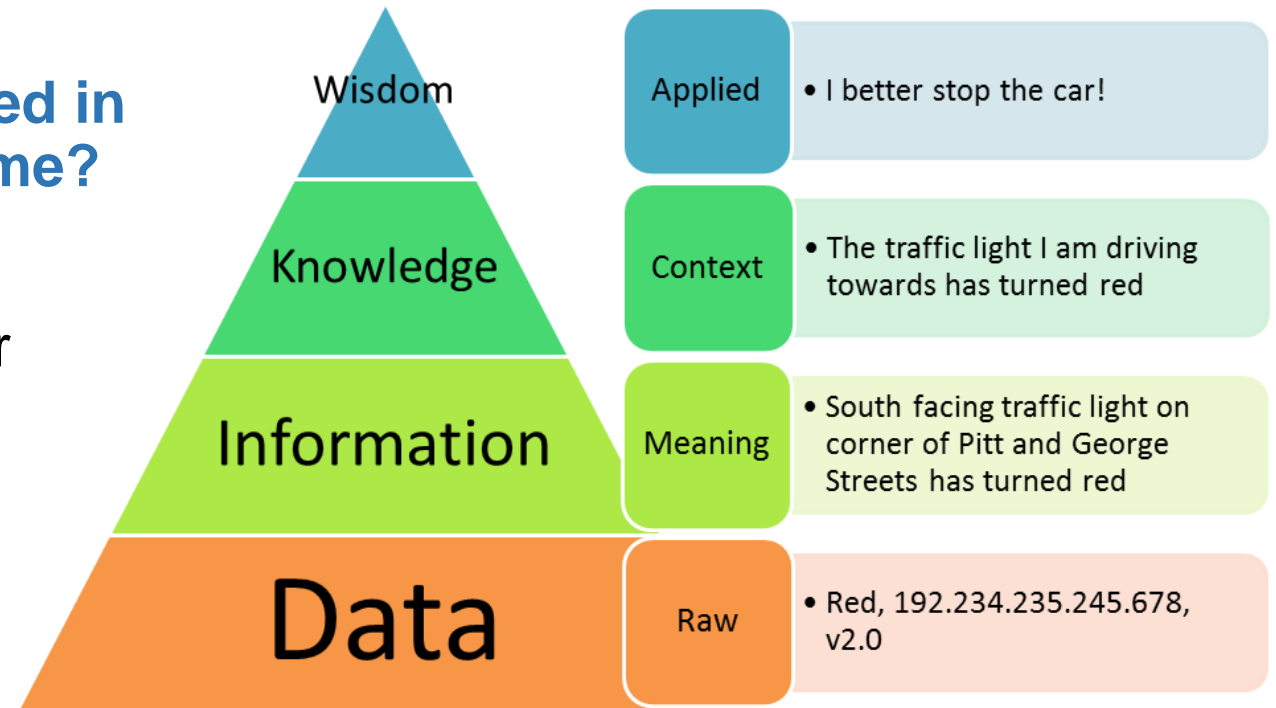
© 2011 Angus McDonald

# Data Use Challenges

To address these challenges, what difficulties are presented?

How can data and resources be linked in order to improve the data spin-up time?

How can we work to educate unfamiliar users?



© 2011 Angus McDonald

# What is a Knowledge Base?

- Think “**Google Search**”.
- Developed by Google in 2012 to enhance the results of its search engine by systematically linking information.
- Aggregates structured and detailed information about a defined topic.
- Enables users to resolve their query without having to navigate and assemble information manually.
- Why not apply it to Earth science and hydrologic data and information?

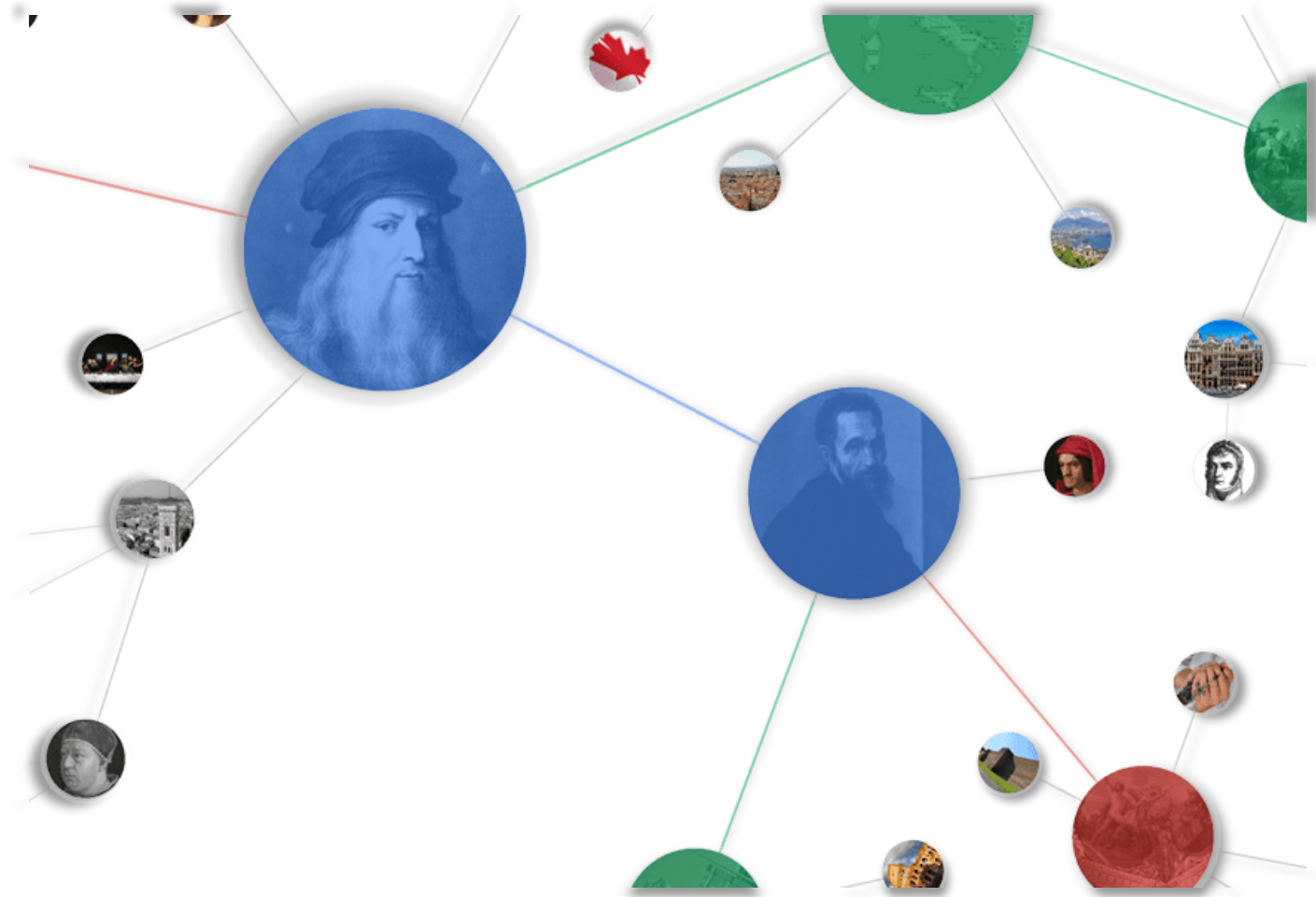


Google Search

I'm Feeling Lucky

# Project Objectives

1. Identify key science information and develop an information model.
2. Extract key information from scientific literature (e.g. hypothesis, conclusions, methods, datasets, variables, etc.).
3. Link scientific knowledge to datasets, resources, services and scientists.
4. Develop a knowledge-based search capability for NASA Earth Science.



[Google, The Knowledge Graph](#)



# Technical Approach

## Terminology

### What is an **entity**?

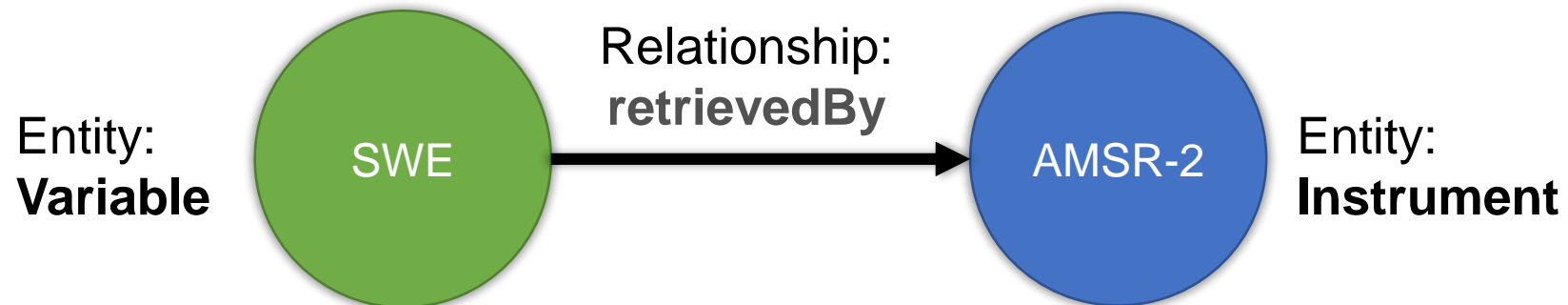
A thing with distinct and independent existence.

Examples: Variables, datasets, instruments, platforms etc.

### What is a **relationship**?

The connection between two entities.

Example: “Snow water equivalent (SWE) is retrieved by AMSR-2”.



# Information Model

The information model defines entities and relationships pertinent to **NASA Earth science and hydrologic data, publications and resources.**





# Technical Approach

## Key Challenge

Knowledge base construction uses both structured and unstructured content (e.g., journal articles).

### Structured Content

- Metadata, tables, controlled vocabularies

```
File "AMSR_2_L3_DailySnow_P00_20170312.he5"
File type: Hierarchical Data Format, version 5

netcdf file:/C:/Users/AWeigel/Downloads/AMSR_2_L3_DailySnow_P00_20170312
{
    group: HDFEOS {

        group: ADDITIONAL {

            group: FILE_ATTRIBUTES {
            }

        }

        group: GRIDS {

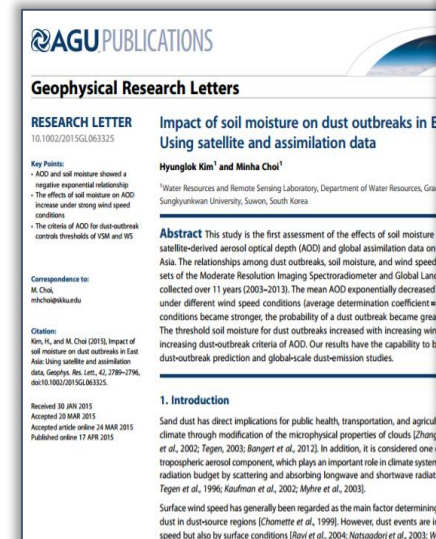
            group: Northern_Hemisphere {
                dimensions:
                    XDim = 721;
                    YDim = 721;
                variables:
                    short _HDFEOS_CRS;
                    :Projection = "HE5_GCTP_LAMAZ";
                    :UpperLeftPointMtrs = -9036843.073845, 9036843.073845; // double
                    :LowerRightMtrs = 9036843.073845, -9036843.073845; // double
                    :ProjParams = 6371228.0, 0.0, 0.0, 0.0, 0.0, 9.0E7, 0.0, 0.0
                    :SphereCode = "19";
            }

        }

        group: Data_Fields {
            variables:
                byte Flags_NorthernDaily(YDim=721, XDim=721);
                :FillValue = -1UB; // byte
                :Unsigned = "true";
            }
        }
    }
}
```

### Unstructured Content

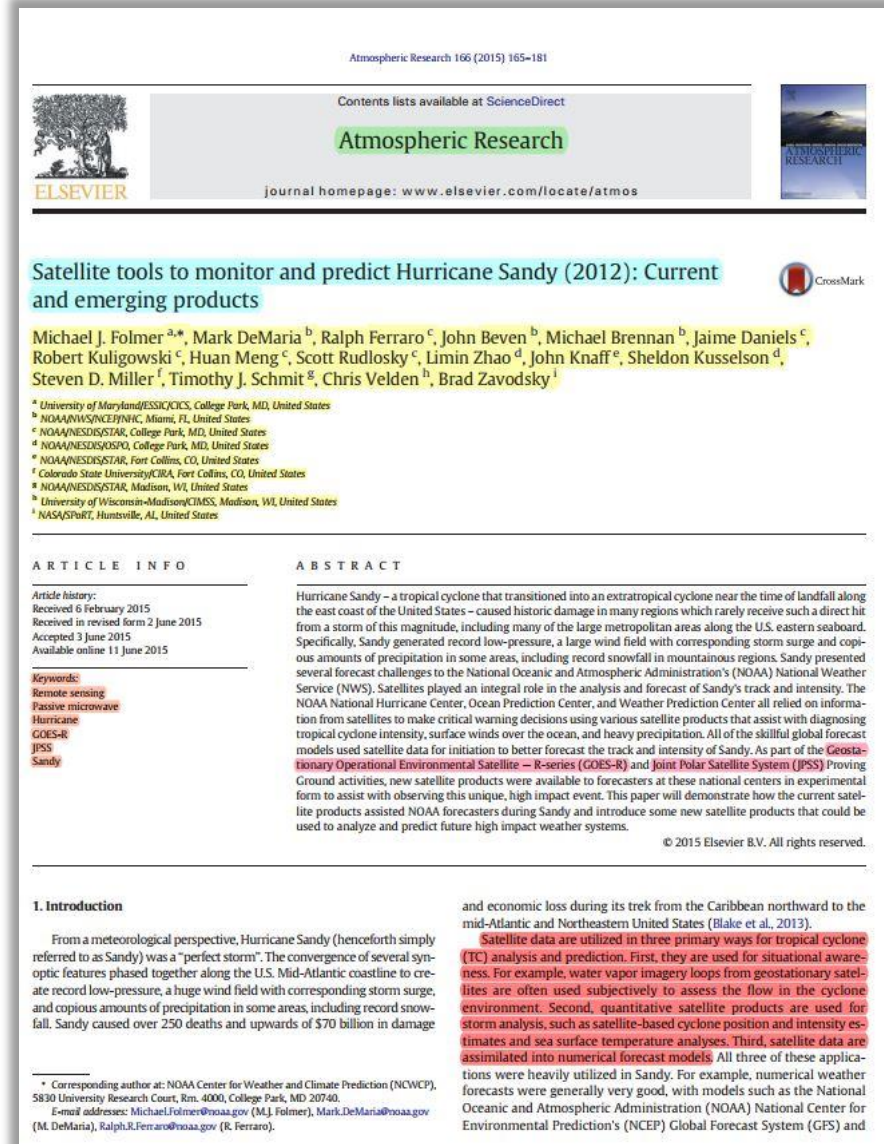
- Journal articles, ATB documents, user guides



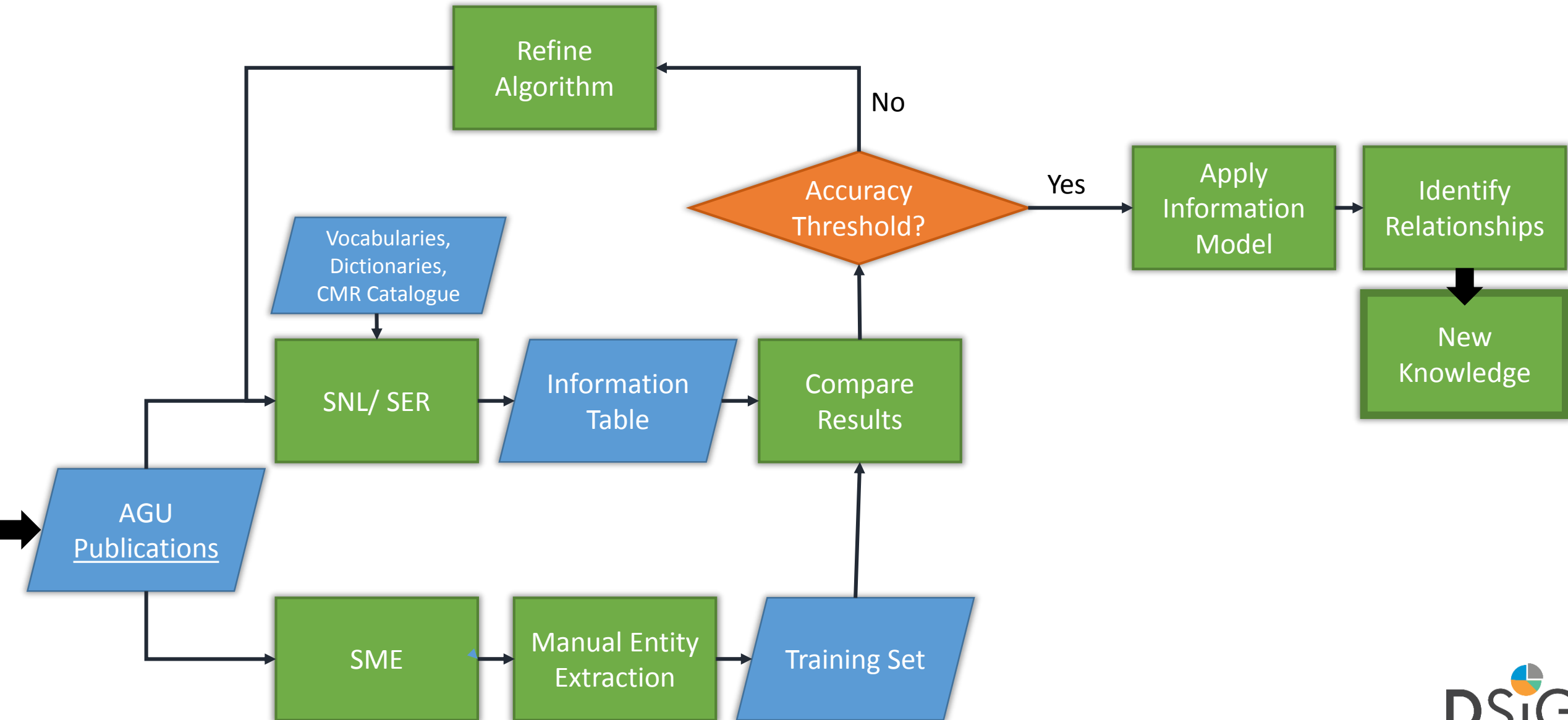
# Technical Approach

## How do we extract entities from unstructured content?

- Natural language processing (**NLP**) – Computers analyzing, understanding and deriving meaning from human language.
- Semantic Entity Recognition (**SER**) – NLP technique used to identify entities in text.
- Use NLP (SER) techniques to identify entities within the unstructured text.
- Apply to journal publication text to extract and identify data, models, methods, people, and institutions (i.e., entities).
- Generate a truth set – Dictionary of known models, science keywords, CMR NASA Earth science data catalogue.




# Technical Approach





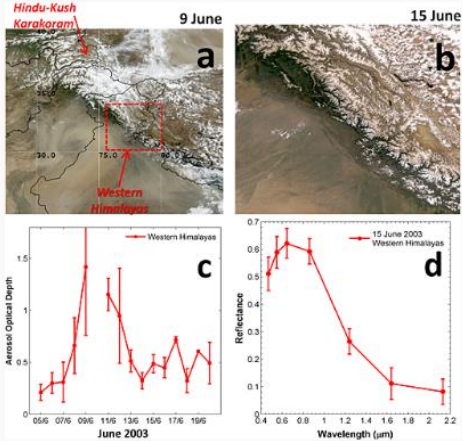
# Preliminary Results

Example of extracted and populated information from unstructured sources

<div><div> Carnegie Mellon University</div><div>Recommendation ▾ Resources ▾ Publication ▾ Analytics ▾ Collaboration ▾ Scientist ▾</div></div>	
Satellite observations of desert dust-induced Himalayan snow darkening	
MetaData	Figures and Captions Analytics Sections Micro Article References Evaluation Proofread
Paper Title	Satellite observations of desert dust-induced Himalayan snow darkening
Authors	Ritesh Gautam , Teppei J. Yasunari , Ritesh Gautam , N. Christina Hsu , William K.-M. Lau , Teppei J. Yasunari
Missing Acronym	BC ( black carbon ); HTP ( Himalaya-Tibetan Plateau ); WH ( western Himalaya ); SCF ( snow cover fraction ); IGP ( Indo-Gangetic Plains ); WH ( western Himalaya ); AOD ( Aerosol Optical Depth ); (NIR) ( n ); TOA ( top-of-atmosphere ); $\mu\text{m}$ ( ); RTM ( radiative transfer model ); SSA ( single scattering albedo );
Id	55
Title	Satellite observations of desert dust-induced Himalayan snow darkening
Year	2013
Date	16-Mar 2013
Url	<a href="http://onlinelibrary.wiley.com/doi/10.1002/grl.50226/full">http://onlinelibrary.wiley.com/doi/10.1002/grl.50226/full</a>
Document Id	55
Flag	0
Subject Category	Geology
Channel Name	GEOPHYSICAL RESEARCH LETTERS
Key Words	HYDROLOGICAL CYCLE; SPECTRAL ALBEDO; TIBETAN PLATEAU; MONSOON; AEROSOLS; PRODUCTS; GLACIERS; IMPACTS; CLIMATE; COVER
Author Keyword	Dust, Snow; Remote Sensing

Topic	climate(87.45%); dust(12.54%); aerosols(0.01%);
Topic Keywords	HYDROLOGICAL CYCLE; SPECTRAL ALBEDO; TIBETAN PLATEAU; MONSOON; AEROSOLS; PRODUCTS; GLACIERS; IMPACTS; CLIMATE; COVER
Instruments	MODIS(Moderate-Resolution Imaging Spectroradiometer)
Datasets	Transit (C1214558371-NOAA_NCEI) Terra 1 MODIS Imagery (C1214598068-SCIOPS) Aqua 1 MODIS Imagery (C1214598104-SCIOPS)
Variables	Toa cloud radiative effect, Solar zenith angle, Snow grain size, Normalized difference vegetation index, Surface snow area fraction, Tropopause instantaneous radiative forcing
Organizations	TOA, SCF, SSA, IGP, HTP
Sponsored Organizations/Projects	NASA

## 3 Extracted Figure/Table

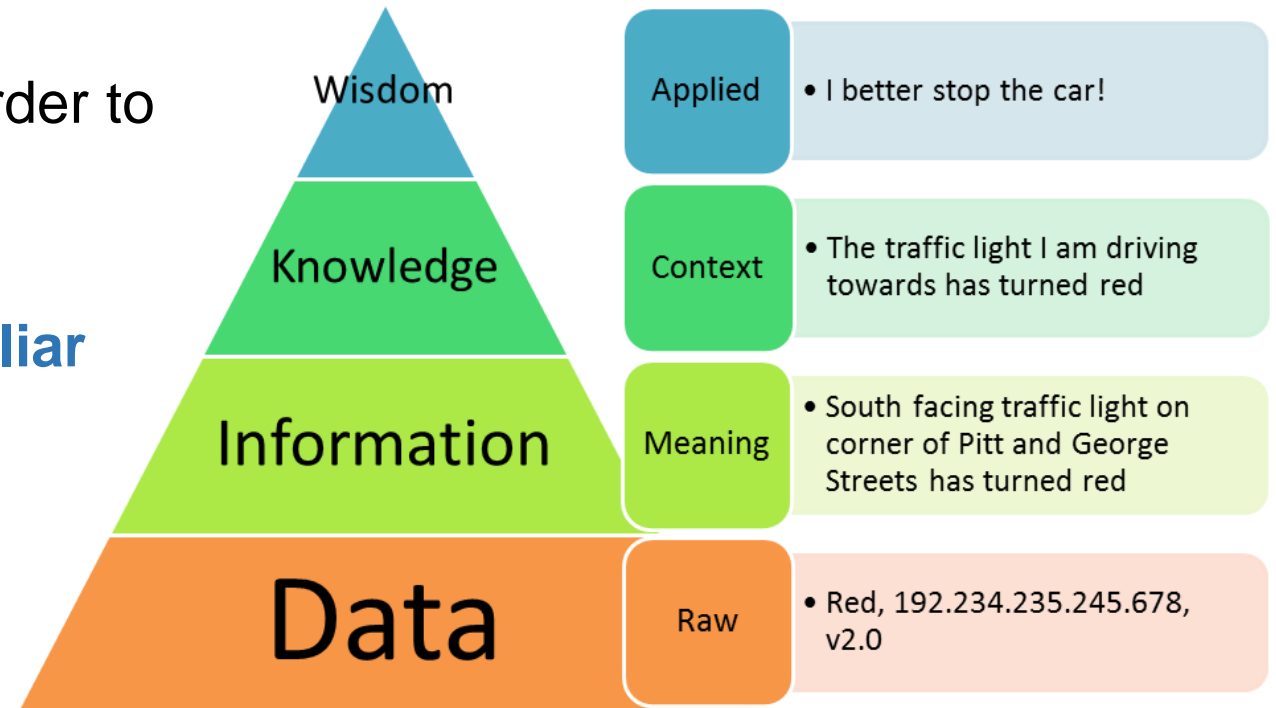
Extracted Figure/Table	Caption
	(a) Satellite image of a major dust outbreak over south Asia, on 9 June 2003 from Terra/MODIS, indicating visibly dust-laden snow surface in the western Himalaya (WH); (b) zoom-in over WH on 15 June 2003; (c) daily AOD variations over the foothills, south of the WH snow cover; (d) MODIS spectral surface reflectance on 15 June indicating the VIS-NIR gradient for WH (30°N–34°N, 76°E–80°E), with error bars of $\pm 1\sigma$ representing pixel-level variability.

# Data Use Challenges

**To address these challenges, what difficulties are presented?**

How can data resources be linked in order to improve the data spin-up time?

How can we work to educate unfamiliar users?



© 2011 Angus McDonald

# Other Resources

Challenge - Publications and technical documents often prove difficult for new and unfamiliar users to digest.

NASA Global Hydrology Resource Center (**GHRC**) Data Active Archive Center (**DAAC**)

NASA Short-term Prediction Research and Transition Center (**SPoRT**)

**What resources are available to introduce data, methods and concepts?**

- GHRC DAAC Data Recipes
- GHRC DAAC Micro Articles
- NASA SPoRT Quick Guides





# NASA GHRC DAAC Data Recipes

## What is a Data Recipe?

Tutorials or step-by-step instructions to help users learn how to discover, visualize and use data, information, software and techniques.

## Types of Data Recipes

- Using netCDF data in ArcGIS
- GHRC tool tutorials
- Python notebooks and scripts
- Data format conversions and georeferencing

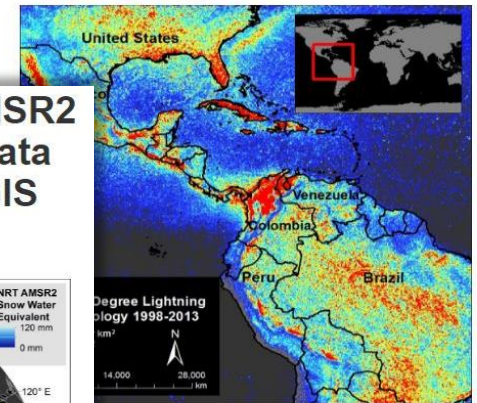
[Discover GHRC Data Recipes Here](#)

## Using ArcGIS to Convert LIS Very High Resolution Gridded Lightning Climatology NetCDF Data to GeoTIFF Format

[Description](#) | [How to Use](#) | [Dataset Information](#) | [Key Parameters](#)

### Description

The Lightning Imaging Sensor (LIS) aboard the Tropical Rainfall



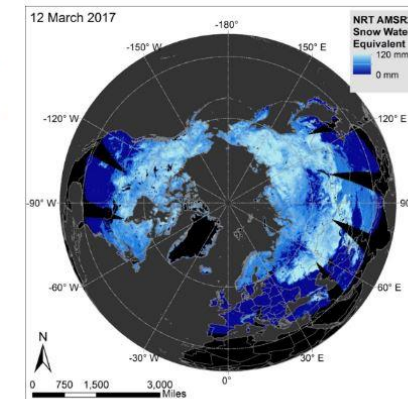
## How to Georeference and Convert NRT AMSR2 Snow Water Equivalent Polar EASE-Grid Data to GeoTIFF Format using Python and ArcGIS

[Description](#) | [How to Use](#) | [Dataset Information](#) | [Key Parameters](#)

### Description

The near real-time (NRT) Land Atmosphere Near real-time Capability for EOS (LANCE) AMSR2 Daily Global Snow Water Equivalent (SWE) EASE-Grids dataset contains SWE and quality assurance flag information for the Northern and Southern Hemispheres. These data are available in HDF-EOS5 format requiring georeferencing in order to display according to the NSIDC EASE-Grid format. This data recipe employs Python to georeference and create four GeoTIFFs of the SWE and quality assessment flag layers within the HDF-EOS5 files. The data are brought into ArcMap where the projection is defined to enable further data processing, analysis and map making. This data recipe requires a pre-installed version of ArcMap, Python and the necessary Python packages.

*Image created using the NRT AMSR2 Snow Water Equivalent Polar EASE-Grid dataset in ArcMap 10.2*



### Data Recipe Type



Data Format Conversion

### Supporting Software Information



TYPE  
Python Script



ArcMap 10.2+



ACCESS  
Open Source



Restricted, license required

# NASA GHRC DAAC Micro Articles

## What is a Micro Article?

A short, interesting document that brings together data and key science concepts

Creates a knowledge base for users by curating around GHRC's data and science thematic areas


## Types of Micro Articles

- Instruments
- Phenomena
- Events or Case Studies
- Publications

[Discover Micro Articles Here](#)


Home >> GHRC Micro Articles >> Phenomenon >> Lightning

## LIGHTNING

 Atmospheric Phenomenon

### WHAT IS LIGHTNING?

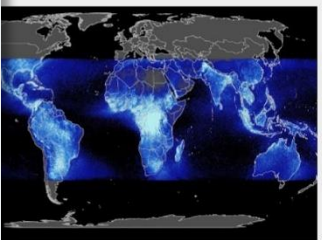
Lightning is the electrical discharge between positively and negatively charged regions within clouds. The electrical discharge serves as an equalization process between the charged regions, and can travel from cloud-to-cloud, cloud-to-ground, or cloud-to-air. Visually, lightning appears as a bright flash of light, or a stroke.




[https://commons.wikimedia.org/wiki/File:Pink\\_Lightning.jpg](https://commons.wikimedia.org/wiki/File:Pink_Lightning.jpg)


### How does lightning occur?

The particles within a cloud interact with each other, causing the particles to fracture and break apart. It is currently believed that smaller ice particles tend to acquire a positive charge, while the larger particles acquire a more negative charge. Under the influences of thunderstorm updrafts, these particles separate until the upper portion of the cloud acquires a net positive charge, and the lower portion of the cloud becomes negatively charged. This separation of charge creates an electrical potential both within the cloud and between the cloud and ground. Eventually, the electrical resistance in the air between the charged regions breaks down and a flash begins. Lightning strokes are an electrical discharge between positive and negative regions of a thunderstorm.



<http://dx.doi.org/10.5067/LIS/LIS/DATA301>


 Volcanic eruptions

 Wildfires

## Lake Effect Snow Event during GCPEX Field Campaign

### Event

**What happened and why it happened**  
A lake effect snow event occurred during the GPM Cold-season Precipitation Experiment (GCPEX) field campaign in Ontario, Canada during the 2011-2012 winter season. Cold, northwest winds moved across the Georgian Bay in eastern Lake Huron and picked up moisture from the lake that fed the development of clouds and snow that varied considerably across a small region south of the bay. The snow clouds developed into persistent narrow bands that resulted in 2 inches of snowfall accumulation at one of the ground sites whereas only 12 miles away they produced 16 inches of snow.



### Science Question

Lake effect snow is generated when cold air moves over warm lake waters such that narrow bands of snow clouds develop. The warmer lake waters heat the lower portions of air causing it to become less dense and begin to rise. As this moisture-laden, warmer air rises it begins to cool leading to condensation and the formation of clouds that can become rather tall enabling the growth of very large snowflakes. Lake effect snow bands can produce snowfall rates exceeding 5 inches an hour, especially if the wind is directed along the largest width of the lake so that a great deal of moisture is continually supplied to the clouds.

### Get Data

The GPM Cold Season Precipitation Experiment (GCPEX) was a field campaign that occurred in Ontario, Canada during the 2012 winter season. The objective of the GCPEX campaign was to study snowfall's physical and radiative properties from the ground through the atmosphere. These measurements are used to help scientists understand the minimum snow rate that can be detected from space and also how well space sensors can discriminate between snow, rain and clear air. Measurements were taken from five ground sites and three research aircraft to provide as complete a sampling as possible.

**SPATIAL COVERAGE**  
[N: 47, W: -80.2, E: -67.7, S: 43.5] degrees

**TIME RANGE**  
February 10-12, 2012

**EVENT TYPE**  
Lake Effect Snow



# NASA SPoRT Quick Guides

## What is a Quick Guide?

Short, easy to use resources that highlight key aspects of a data product or tool.

Intended to assist forecasters in quickly recalling information during times of operation.

## Available Forms

- Download/print
- Interactive web browser
- Interactive through personal display system

[Discover SPoRT Quick Guides here](#)

**True Color RGB**

**Daytime Convection RGB**

**Quick Guide**

**Why is the Convection RGB imagery Important?**

The Daytime Convective Storms (Convection) RGB was designed for identification of convection with strong updrafts and small ice particles indicative of severe storms. This RGB helps increase nowcasting capabilities of severe storms by identifying the early stage of strong convection. Knowing the microphysical characteristics of convective clouds helps determine storm strength and stage to improve nowcasts and short-term forecasts. Bright yellow in the RGB indicates strong updrafts prior to the mature storm stage.

**Convection RGB Recipe**

Color	Band / Band Diff. (μm)	Physically Relates to...	Small contribution to pixel indicates...	Large Contribution to pixel indicates...
Red	6.2 – 7.3	Cloud height	Low clouds	High clouds
Green	3.9 – 10.4	Particle size	Large ice or water particles, weak updrafts	Small ice or water particles, strong updrafts
Blue	1.6 – 0.64	Cloud phase	Ice clouds	Water clouds

**Impact on Operations**

**Primary Application**  
**Convection and Severe Weather:** identify intense updrafts that indicate strong convection.  
Strong convection is bright yellow.

**Limitations**

**Daytime only application:** the RGB relies on solar reflectance from visible, near-IR, and shortwave IR channels. Day color impacted by sun/satellite viewing.

**Limitations**

Small contribution to pixel indicates ...  
Large contribution to pixel indicates ...  
Low reflectance for the given band  
High reflectance for the given band

**Limitations**

True color reflectance channels, and available at night.

**Maps:** China, Tianjin, Yellow Sea, Residual Smoke (Explosion), Mongolia, Kazakhstan, China, Dust Storm.

Himawari-8 AHI (courtesy CIRA) showing (a) convection, (b) biomass smoke, and (d) lofted dust.



# Other Resources

## **How do these benefit Earth science and hydrologic applications?**

- These resources provide introductory information that is easy to read and understand without overwhelming users.
- Each point to additional documents for more detailed information.
- Each contain information on commonly used data, models, and software.
- They link directly to data, helping users understand a dataset and how to apply it towards research or applications.
- Populating this information within a knowledge graph allows users to search and discover information on data and methods for a broad user community.

# Next Steps

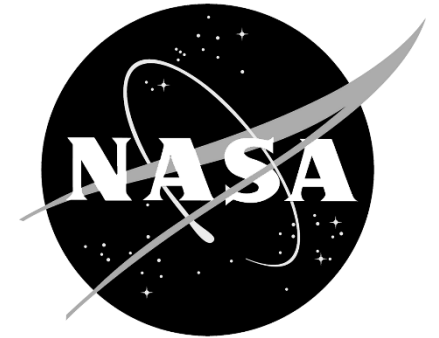
- Investigate generating easy to understand resources in a structured format to allow more seamless integration within the knowledge graph.
- Continue refining SER for Earth Science
- Continue building and evaluate a training set for SER (working with graduate students and SMEs)
- Scale efforts to to all Earth Science related journal titles in the Wiley Online Library
- Begin mining graphs to obtain new information
  - Prediction of relationship between entities (i.e., Network Link Prediction)
  - Automatic generation of new content (e.g., MicroArticles)

# Benefits to NASA Earth Science and Hydrology

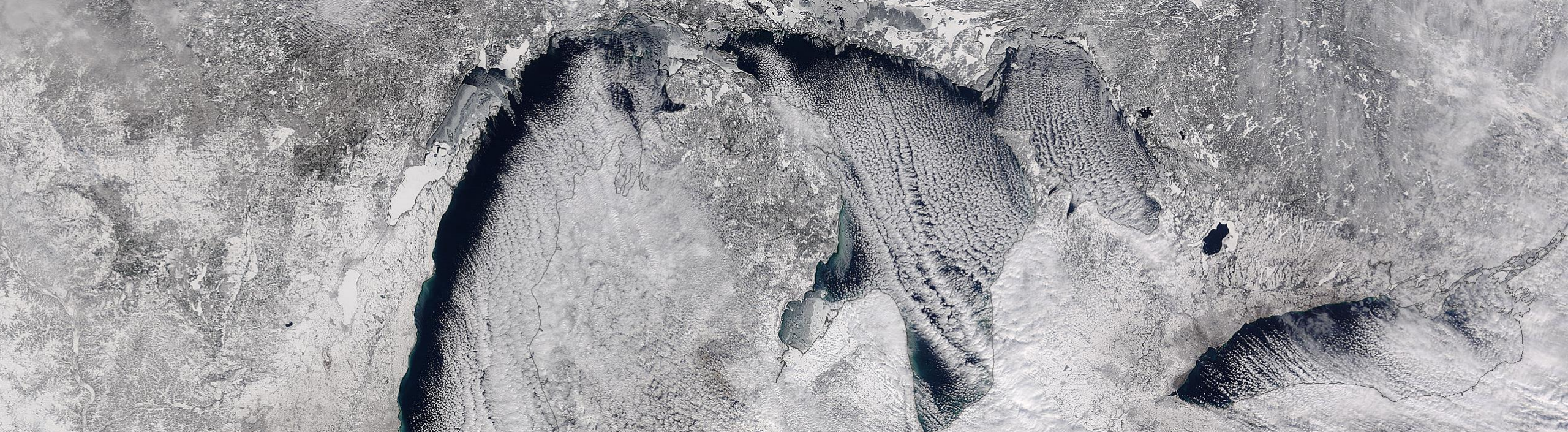
- Addresses the challenge in navigating the increasing volume of data and information.
- To provide an operational knowledge base to enhance NASA's Earth science research.

## Beneficial Applications

- Hypothesis formulation and testing:
  - Automate the search for and compilation of background information.
  - Given a topic, what hypotheses have been tested?
  - What data/tools are being used to test a hypothesis?
  - Common paths to knowledge discovery.
- Mission development/review:
  - What kinds of instruments/parameters are needed to specify science objectives?
  - Impact of a mission by linking it with publications and dataset distribution.







# Thank you, questions?

[amanda.m.weigel@nasa.gov](mailto:amanda.m.weigel@nasa.gov)

NASA/MSFC Data Science Informatics Group

